# EXPLORING LEARNABILITY IN MEMORY-AUGMENTED RECURRENT NEURAL NETWORKS: PRECISION, STABILITY, AND EMPIRICAL INSIGHTS

**Shrabon Das**
University of South Florida
Tampa, FL 33620, USA
{das157}@usf.edu

**Ankur Mali**
University of South Florida
Tampa, FL 33620, USA
{ankurarjunmali}@usf.edu

## ABSTRACT

This study investigates the learnability of memory-less and memory-augmented Recurrent Neural Networks (RNNs) with deterministic and non-deterministic stacks, which are theoretically equivalent to Pushdown Automata in terms of expressivity. However, empirical evaluations reveal that these models often fail to generalize on longer sequences, particularly when learning context-sensitive languages, suggesting they rely on precision rather than mastering symbolic grammar rules. Our experiments examined fully trained models and models with various frozen components: the controller, the memory, and only the classification layer. While all models showed similar performance on training validation, the model with frozen memory achieved state-of-the-art performance on the Penn Treebank (PTB) dataset, reducing the best overall test perplexity from 123.5 to 120.5—a gain of approximately 1.73%. When tested on context-sensitive languages, models with frozen memory consistently outperformed others on small to medium test sets. Notably, well-trained models experienced up to a 60% performance drop on longer sequences, whereas models with frozen memory retained close to 90% of their initial performance. Theoretically, we explain that freezing the memory component enhances stability by anchoring the model's capacity to manage temporal dependencies without constantly adjusting memory states. This approach allows the model to focus on refining other components, leading to more robust convergence to optimal solutions. These findings highlight the importance of designing stable memory architectures and underscore the need to evaluate models on longer sequences to truly understand their learnability behavior and limitations. The study suggests that RNNs may rely more on precision in data processing than on internalizing grammatical rules, emphasizing the need for improvements in model architecture and evaluation methods.

## Introduction

Recurrent Neural Networks (RNNs) have been foundational in sequence modeling due to their ability to capture temporal dependencies. Architectures such as Elman RNNs, Gated Recurrent Units (GRUs), and Long Short-Term Memory networks (LSTMs) [1] are widely used in applications like speech recognition, machine translation, and time-series analysis. However, these models are constrained by their fixed memory capacity, limiting them to recognizing regular languages when implemented with finite precision [2, 3].

To enhance the computational capabilities of RNNs, researchers have explored augmenting them with external memory structures like stacks [4, 5, 6, 7, 8, 9, 10]. This approach extends the expressivity of RNNs to context-free languages (CFLs) [11], which are crucial in applications like natural language processing (NLP) where hierarchical structures are prevalent.

Memory-augmented models have demonstrated significant improvements in recognizing complex formal languages by simulating operations similar to Pushdown Automata (PDA). These models can process deterministic and nondeterministic context-free languages, thereby expanding the class of languages they can handle compared to traditional RNNs. However, while their theoretical expressivity is well-established, the learnability of these sys-

tems remains an open challenge. Learnability here refers to the model's ability to reliably generalize from training data to unseen inputs, particularly with complex recursive patterns.

Despite their theoretical advantages, empirical studies indicate that many memory-augmented models struggle to generalize to longer sequences. This instability is likely due to the dynamic nature of memory manipulation and finite precision constraints, leading to performance degradation on extended sequences. Learnability is closely tied to stability: unstable systems are prone to accumulating errors over time, particularly as sequence lengths increase, resulting in unpredictable behavior. Understanding stability conditions is therefore crucial for enhancing learnability. In this work, we investigate how different configurations, such as freezing the RNN controller while training the memory, impact stability and performance compared to fully trainable but unstable setups. Stability, in this context, refers to a model's ability to maintain consistent performance across varying sequence lengths and input complexities. A stable model efficiently generalizes from training data while remaining robust to variations in input and computational precision. Stability is key to ensuring reliable learnability for complex tasks.

Key aspects of stability include:

- **Consistency Across Sequence Lengths:** Ensuring stable performance on both short and long sequences.
- **Robustness to Variability:** Handling variations in input and precision without significant performance degradation.
- **Reliable Memory Manipulation:** Maintaining correct memory operations even under finite precision.
- **Enhanced Generalization:** Promoting better generalization across different tasks and datasets.

This study addresses the critical relationship between stability and learnability in stack-augmented RNNs by providing theoretical and empirical insights into when these models succeed and when they fail. Our key contributions include:

- Theoretical analysis of stability and instability conditions for stack-augmented RNNs, showing that configurations with a frozen RNN controller but a trainable stack can outperform fully trained yet unstable models.
- Derivation of error bounds for unstable systems, illustrating how models that initially perform well on slightly longer sequences may eventually converge to random guessing as sequence lengths increase.
- Analysis of learnability under varying conditions, demonstrating how stability plays a crucial role in effective generalization to unseen sequences.
- A framework for understanding the impact of machine precision on memory operations, emphasizing the need for stable memory manipulation to ensure robustness.
- Empirical validation through experiments on context-sensitive languages and the Penn Treebank dataset, demonstrating alignment with theoretical findings.

By focusing on stability, we bridge the gap between theoretical expressivity and practical learnability, offering insights into designing more robust architectures capable of handling complex language structures while maintaining stable performance across diverse conditions. Understanding stability through the lens of machine precision and error growth provides strategies for improving the learnability of stack-augmented RNNs in real-world applications

## Background

This section describes the basic concepts discussed in this work, focusing on Pushdown Automata (PDA) and stack-augmented Recurrent Neural Networks (RNNs). We explore how these models differ in expressivity and how they recognize formal languages. **Pushdown Automaton (PDA):** A PDA is a computational model capable of recognizing context-free languages using a stack as auxiliary memory. Formally, a PDA is defined as:

$$M = (Q, \Sigma, \Gamma, \delta, q_0, Z_0, F),$$

where $Q$ represents states, $\Sigma$ the input alphabet, $\Gamma$ the stack alphabet, and $\delta$ the transition function. The PDA reads an input string $x \in \Sigma^*$, updates its stack and states according to $\delta$, and either accepts or rejects $x$ based on the final state and stack configuration.

Stack-Augmented Recurrent Neural Network (RNN with Memory): To extend RNNs' expressivity to context-free languages, they can be augmented with a stack. A stack-augmented RNN is defined as:

$$f_\theta = (\theta_c, \theta_m),$$

where $\theta_c$ governs the RNN controller, and $\theta_m$ controls stack operations (push, pop, no-op). At each time step, the model:

1. Updates the hidden state:
$$h_t = f(W_h h_{t-1} + W_x x_t + b).$$

2. Chooses a stack operation:
$$\text{Stack}_{t+1} = \delta_s(h_t, \text{Stack}_t),$$
where $\delta_s$ is a learned function.

3. Produces the output based on both the hidden state and stack:
$$y_t = g(W_y h_t + W_s \text{Stack}_t + b_y).$$

Comparison of Expressivity: While a PDA directly handles context-free languages using a stack, a stack-augmented RNN approximates this behavior through learned stack operations. The primary challenge lies in maintaining stability and generalization as sequence complexity increases, particularly for long or deeply nested structures.

## Methodology

In this section, we first formally define the stability condition of memory less model, where stability is defined as follows

**Definition 0.1** (Stability). Let $L$ be a formal language accepted by a discrete state machine $M$, such as a finite automaton. A sequence model $f$ is said to be *stable* with respect to $L$ if, for any input sequence $x = (x_1, x_2, \ldots, x_T)$ of arbitrary length $T \in \mathbb{N}$:

- There exists a mapping $\phi : \text{States}(M) \to \text{States}(f)$ such that the state transitions of $f$ are functionally equivalent to those of $M$. Formally, for each state $q \in \text{States}(M)$ and corresponding state $\phi(q) \in \text{States}(f)$:
$$\delta_M(q, x_i) = q' \iff \delta_f(\phi(q), x_i) = \phi(q'),$$
where $\delta_M$ and $\delta_f$ denote the transition functions for $M$ and $f$, respectively, and $x_i$ is an input symbol.

- The acceptance condition for $L$ is preserved by $f$. That is, for any sequence $x \in \Sigma^*$:
$$x \in L \iff f(x) = 1,$$
where $\Sigma$ is the input alphabet, and $f(x) = 1$ indicates acceptance by the model $f$.

- The behavior of $f$ remains consistent for sequences of arbitrary length. For any two sequences $x^{(1)}$ and $x^{(2)}$ of lengths $T_1$ and $T_2$, respectively, with $T_1 \neq T_2$:
$$M(x^{(1)}) = f(x^{(1)}) \quad \text{and} \quad M(x^{(2)}) = f(x^{(2)}),$$
where $M(x)$ and $f(x)$ denote the outputs of the discrete state machine and the model, respectively.

- Stability is maintained as $T \to \infty$, ensuring that $f$ does not exhibit performance degradation or erratic behavior for long sequences.

In essence, a model $f$ is stable if it is functionally equivalent to the discrete state machine $M$ that defines the language $L$, preserving both the state transitions and acceptance conditions across sequences of arbitrary length.

Similarly stability of memory augmented RNN is formally defined as follows:

**Definition 0.2** (Stability of a Pushdown Automaton (PDA) Model). Let $L$ be a formal language accepted by a Pushdown Automaton (PDA) $M = (Q, \Sigma, \Gamma, \delta, q_0, Z_0, F)$, where:

- $Q$ is the set of states,

- $\Sigma$ is the input alphabet,

- $\Gamma$ is the stack alphabet,

- $\delta : Q \times (\Sigma \cup \{\epsilon\}) \times \Gamma \to Q \times \Gamma^*$ is the transition function,

- $q_0 \in Q$ is the initial state,

- $Z_0 \in \Gamma$ is the initial stack symbol,

- $F \subseteq Q$ is the set of accepting states.

A sequence model $f$ is said to be *stable* with respect to $L$ if, for any input sequence $x = (x_1, x_2, \ldots, x_T)$ of arbitrary length $T \in \mathbb{N}$:

- **State Stability:** There exists a mapping $\phi : Q \to \text{States}(f)$ such that the state transitions of $f$ are functionally equivalent to those of the PDA $M$. Formally, for each state $q \in Q$ and corresponding state $\phi(q) \in \text{States}(f)$:
$$\delta(q, x_i, \gamma) = (q', \gamma') \iff \delta_f(\phi(q), x_i, \gamma) = (\phi(q'), \gamma'),$$
where $\gamma \in \Gamma$ is the top symbol on the stack and $\delta_f$ represents the transition function of the model $f$.

- **Stack Stability:** For the stack operations in $f$ to be stable, the stack manipulation should be equivalent to that of the PDA $M$. Let $\text{Stack}_M(t)$ and $\text{Stack}_f(t)$ represent the stack contents at time $t$ for $M$ and $f$, respectively. The model $f$ is stable if:
$$\text{Stack}_M(t) = \text{Stack}_f(t) \quad \text{for all } t.$$

This ensures that push and pop operations are executed consistently and that the stack remains in sync with the PDA for arbitrary input sequences.

- **Acceptance Condition:** The model $f$ correctly accepts or rejects sequences according to $L$. That is, for any sequence $x \in \Sigma^*$:
$$x \in L \iff f(x) = 1,$$
where $f(x) = 1$ indicates acceptance by the model $f$.

- **Behavioral Consistency Across Sequence Lengths:** For any two sequences $x^{(1)}$ and $x^{(2)}$ of lengths $T_1$ and $T_2$, respectively, with $T_1 \neq T_2$:
$$M(x^{(1)}) = f(x^{(1)}) \quad \text{and} \quad M(x^{(2)}) = f(x^{(2)}),$$
where $M(x)$ and $f(x)$ denote the outputs of the PDA and the model, respectively.

- **Long-Term Stability:** The model $f$ remains stable for arbitrarily long sequences, ensuring that its behavior and stack operations do not degrade or become inconsistent as $T \to \infty$.

In summary, a model $f$ is stable if it is functionally equivalent to the PDA $M$ that defines the language $L$, preserving both the state transitions and stack operations across sequences of arbitrary length.

**Theoretical Analysis of Stability of Memory Augmented Neural Network**

In this section, we describe the analytical approach used to understand the behavior of stack-augmented Recurrent Neural Networks (RNNs) when processing formal languages. Our analysis is structured around a series of theorems that establish key properties such as error growth, stability, and the comparative performance of different configurations of stack-augmented RNNs.

First, we formalize the relationship between stability and model performance by evaluating the impact of stack operations and state transitions. Stability is crucial for ensuring consistent performance across sequences of varying lengths, and the presented theorems highlight how instability leads to unbounded error growth, ultimately converging to random guessing or worse. Formally we can show the existence of stable memory augmented RNN as follows:

**Theorem 0.3** (Stability of Stack-Augmented RNNs). *Let $L$ be a formal language recognized by a Pushdown Automaton (PDA) $M = (Q, \Sigma, \Gamma, \delta, q_0, Z_0, F)$. Consider a Recurrent Neural Network (RNN) $f$ augmented with a stack that models the PDA $M$. The RNN $f$ is said to be stable if it satisfies the following conditions:*

1. *The state transitions of $f$ are functionally equivalent to those of the PDA $M$.*

2. *The stack operations of $f$ (push, pop, no-op) are consistent with those of $M$.*

3. *The outputs of $f$ and $M$ are equivalent for sequences of different lengths.*

*4. As the sequence length $T \to \infty$, the error between $f$ and $M$ approaches zero.*

*The RNN $f$ is unstable if any of the conditions above are violated.*

*Proof.* 1. Let $\phi : Q \to \text{States}(f)$ be a mapping between the states of the PDA $M$ and the states of the RNN $f$. The state transition function of $M$ is defined as:

$$\delta : Q \times (\Sigma \cup \{\epsilon\}) \times \Gamma \to Q \times \Gamma^*,$$

where $\Sigma$ is the input alphabet, $\Gamma$ is the stack alphabet, and $\delta$ is the transition function of $M$. The transition function of $f$, denoted as $\delta_f$, should satisfy:

$$\delta(q, x_i, \gamma) = (q', \gamma') \iff \delta_f(\phi(q), x_i, \gamma) = (\phi(q'), \gamma'),$$

for all states $q \in Q$, input symbols $x_i \in \Sigma$, and stack symbols $\gamma \in \Gamma$. If such a mapping $\phi$ exists and is consistent across all transitions, then the state transitions of $f$ are stable. The existence of this mapping ensures that for any input sequence, the behavior of $f$ mirrors the PDA $M$, maintaining state stability.

2. The stack operations in $f$ must replicate those of $M$ across all time steps. For any input sequence $x = (x_1, \ldots, x_T)$ and corresponding stack operations $\gamma \in \Gamma$, let $\text{Stack}_M(t)$ and $\text{Stack}_f(t)$ denote the stack contents at time $t$ for $M$ and $f$, respectively. Stack stability requires:

$$\text{Stack}_M(t) = \text{Stack}_f(t) \quad \forall t \in \{1, 2, \ldots, T\}.$$

If this condition holds, then the push, pop, and no-op operations of $f$ are consistent with those of $M$. Inconsistent stack behavior leads to errors in sequence recognition, causing instability.

3. Consider two input sequences $x^{(1)}$ and $x^{(2)}$ with lengths $T_1 \neq T_2$. Stability requires that the outputs of $f$ and $M$ are functionally equivalent:

$$M(x^{(1)}) = f(x^{(1)}) \quad \text{and} \quad M(x^{(2)}) = f(x^{(2)}).$$

This behavioral consistency across sequence lengths follows from stable state transitions and stack operations, ensuring that $f$ correctly processes sequences of varying lengths. If $f$ exhibits degradation in performance for longer sequences, the model is unstable.

4. As sequence length $T \to \infty$, stability requires:

$$\lim_{T \to \infty} \text{Error}(f, M) = 0,$$

where $\text{Error}(f, M)$ measures the divergence between the outputs of $f$ and $M$. If this condition holds, the long-term stability of $f$ is guaranteed. Any error accumulation as $T$ increases would indicate instability in either state transitions or stack operations, leading to performance degradation.

Thus the RNN $f$ is stable if all four conditions are met. If any of these conditions fail, the model becomes unstable, resulting in incorrect behavior or output divergence from the PDA $M$. $\square$

Next, we formally prove that unstable memory-augmented RNN, after arbitrary steps, will eventually become equivalent to a random network, hampering its generalization. Formally we prove following theorem

**Theorem 0.4.** *Let $L$ be a formal language recognized by a Pushdown Automaton (PDA) $M = (Q, \Sigma, \Gamma, \delta, q_0, Z_0, F)$. Consider a Recurrent Neural Network (RNN) $f_\theta$ augmented with a stack and parameterized by $\theta$. Let $f_{random}$ be a randomly initialized network with parameters $\theta_{random}$. If the RNN $f_\theta$ is unstable, then the expected loss of $f_\theta$ is equivalent to the expected loss of $f_{random}$, i.e.,*

$$\mathbb{E}\left[Loss(f_\theta(x), y)\right] \approx \mathbb{E}\left[Loss(f_{random}(x), y)\right],$$

*where $x$ is the input, $y$ is the target output, and $Loss(\cdot, \cdot)$ is a suitable loss function.*

*Proof.* Let $M(x)$ represent the correct output for input $x$ according to the PDA $M$, and let $f_\theta(x)$ represent the output of the RNN $f_\theta$ for the same input.

1. **Instability and Divergence from True Dynamics**:

For an unstable RNN $f_\theta$, instability implies that the state transitions and stack operations diverge from those of the PDA $M$. Specifically, for sufficiently large input sequences $x$ with length $T$, the output of $f_\theta(x)$ diverges from the correct output $M(x)$:

$$\lim_{T \to \infty} |f_\theta(x) - M(x)| \to \infty.$$

The divergence grows as the sequence length increases, leading to increasingly erroneous outputs. In such cases, $f_\theta$ fails to learn the correct language $L$ and instead exhibits behavior that is statistically uncorrelated with $M(x)$.

2. **Equivalence to Random Network Behavior**:

For a randomly initialized network $f_{\text{random}}$, the output behavior is not correlated with the target output $M(x)$. The expected loss for a random network is:

$$\mathbb{E}\left[\text{Loss}(f_{\text{random}}(x), y)\right] = \int_{x \in \Sigma^*} \text{Loss}(f_{\text{random}}(x), y) P(x)\, dx$$
$$\approx \text{Constant}.$$

where $P(x)$ is the distribution over the input space $\Sigma^*$.

Since the unstable network $f_\theta$ also diverges from the correct output, its performance becomes uncorrelated with the target output $y$, resulting in an expected loss similar to that of a random network:

$$\mathbb{E}\left[\text{Loss}(f_\theta(x), y)\right] \approx \text{Constant}.$$

3. **Convergence of Expected Losses**:

As the instability of $f_\theta$ increases, the distribution of its output becomes increasingly similar to that of $f_{\text{random}}$. In the limit, the expected loss of $f_\theta$ approaches the expected loss of $f_{\text{random}}$:

$$\lim_{\text{instability} \to \infty} \mathbb{E}\left[\text{Loss}(f_\theta(x), y)\right] = \mathbb{E}\left[\text{Loss}(f_{\text{random}}(x), y)\right].$$

This indicates that the unstable network is no more learnable than a random network.

Thus the unstable RNN $f_\theta$ behaves like a random network in terms of expected loss, it lacks the capacity to learn the correct language $L$. Thus, the learnability of the unstable system is equivalent to that of a random network. $\quad\square$

Next, we provide error bounds for fully-trained, partially-trained, and frozen networks, which is crucial to derive learnability error bounds in practice, as ideally, the error for the stable system should be much lower compared to these variants. Formally, we show

**Theorem 0.5** (Learnability of Gradient-Descent Trained Systems). *Let $L$ be a formal language recognized by a Pushdown Automaton (PDA) $M = (Q, \Sigma, \Gamma, \delta, q_0, Z_0, F)$. Consider a stack-augmented Recurrent Neural Network (RNN) $f_\theta$ parameterized by $\theta = (\theta_c, \theta_m)$, where $\theta_c$ represents the parameters of the RNN/controller and $\theta_m$ represents the parameters governing the stack/memory operations.*

*Define the following configurations:*

- $f_\theta^{full}$*: Both controller $\theta_c$ and memory $\theta_m$ are trainable.*

- $f_{\theta_c}^{frozen}$*: The controller $\theta_c$ is frozen, and only the memory $\theta_m$ is trainable.*

- $f_{\theta_m}^{frozen}$*: The memory $\theta_m$ is frozen, and only the controller $\theta_c$ is trainable.*

- $f_\theta^{frozen}$*: Both the controller $\theta_c$ and the memory $\theta_m$ are frozen (i.e., the system is untrained or random).*

*Let the loss function be $\text{Loss}(f_\theta(x), y)$, where $x \in \Sigma^*$ is an input sequence, and $y$ is the target output. The following conditions hold:*

*1. Stable System Minimizes Error:*

*A stable system, where both $\theta_c$ and $\theta_m$ are trained, achieves the lowest expected loss:*

$$\mathbb{E}\left[Loss(f_\theta^{full}(x), y)\right] < \min\left\{\mathbb{E}\left[Loss(f_{\theta_c}^{frozen}(x), y)\right],\right.$$
$$\mathbb{E}\left[Loss(f_{\theta_m}^{frozen}(x), y)\right],$$
$$\left.\mathbb{E}\left[Loss(f_\theta^{frozen}(x), y)\right]\right\}.$$

*2. **Unstable System's Error is Similar to Untrained System**:*

*If the system is unstable, its behavior resembles that of an untrained or random network, leading to an expected error similar to a completely frozen system:*

$$\mathbb{E}\left[Loss(f_\theta^{unstable}(x), y)\right] \approx \mathbb{E}\left[Loss(f_\theta^{frozen}(x), y)\right].$$

*3. **Bounded Error Growth for Stable Systems**:*

*A stable system exhibits bounded error growth as sequence length increases. There exists a constant $C > 0$ such that:*

$$\left|Loss(f_\theta^{full}(x), y)\right| \leq C \quad \forall T \in \mathbb{N}.$$

*4. **Consistency Across Sequence Lengths**:*

*A stable system exhibits consistent error across sequences of varying lengths:*

$$Var\left(Loss(f_\theta^{full}(x), y)\right) < Var\left(Loss(f_{\theta_c}^{frozen}(x), y)\right) \quad \forall T \in \mathbb{N},$$

*where $Var(\cdot)$ represents the variance of the error as a function of sequence length.*

*Proof.* 1. **A Stable System Minimizes Error**:

In a stable system, gradient descent optimizes both $\theta_c$ and $\theta_m$ to minimize the loss function:

$$\theta^* = \arg\min_\theta \mathbb{E}\left[Loss(f_\theta(x), y)\right].$$

Since both $\theta_c$ and $\theta_m$ are jointly optimized, the fully trained system $f_\theta^{\text{full}}$ achieves lower loss than any system where one or both components are frozen:

$$\mathbb{E}\left[Loss(f_\theta^{\text{full}}(x), y)\right] < \min\left\{\mathbb{E}\left[Loss(f_{\theta_c}^{\text{frozen}}(x), y)\right],\right.$$
$$\mathbb{E}\left[Loss(f_{\theta_m}^{\text{frozen}}(x), y)\right],$$
$$\left.\mathbb{E}\left[Loss(f_\theta^{\text{frozen}}(x), y)\right]\right\}.$$

2. **Next we show that Unstable System's Error is Similar to Untrained System**:

If the system is unstable, it fails to correctly model the state transitions and stack operations. Let the target output be $M(x)$, where $M$ represents the PDA. For an unstable system, the output diverges from $M(x)$ as the sequence length $T$ increases:

$$\lim_{T \to \infty} |f_\theta(x) - M(x)| \to \infty.$$

Consequently, the system behaves similarly to a random or untrained network, leading to:

$$\mathbb{E}\left[Loss(f_\theta^{\text{unstable}}(x), y)\right] \approx \mathbb{E}\left[Loss(f_\theta^{\text{frozen}}(x), y)\right].$$

3. **Next we show that Bounded Error Growth for Stable Systems**:

For a stable system, the optimization process ensures that the error remains bounded across varying sequence lengths. There exists a constant $C > 0$ such that:

$$\left| \text{Loss}(f_\theta^{\text{full}}(x), y) \right| \leq C \quad \forall T \in \mathbb{N}.$$

This condition is satisfied because the model correctly captures the underlying dynamics of the PDA, preventing unbounded error growth.

4. **Next we show the the Consistency Across Sequence Lengths**:

In a stable system, both the controller and memory are optimized to handle sequences of varying lengths. This results in low variance in the error:

$$\text{Var}\left(\text{Loss}(f_\theta^{\text{full}}(x), y)\right) < \text{Var}\left(\text{Loss}(f_{\theta_c}^{\text{frozen}}(x), y)\right).$$

The low variance implies that the system is robust to changes in sequence length, a key characteristic of stability.

$\square$

Next, we explore the conditions under which a frozen RNN with a trainable stack can outperform an unstable fully trained model. This comparison is significant for understanding how different components (controller versus memory) contribute to overall performance, especially when the model faces complex or recursive patterns inherent in context-free languages. The analysis shows that freezing the RNN while allowing the stack to adapt can stabilize performance, leading to more reliable error bounds.

**Theorem 0.6** (Frozen RNN with Trainable Memory Outperforms an Unstable Fully-Trained Model). *Let $L$ be a formal language recognized by a Pushdown Automaton (PDA) $M = (Q, \Sigma, \Gamma, \delta, q_0, Z_0, F)$. Consider a stack-augmented Recurrent Neural Network (RNN) $f_\theta$ parameterized by $\theta = (\theta_c, \theta_m)$, where:*

- *$\theta_c$ represents the parameters of the RNN/controller.*

- *$\theta_m$ represents the parameters governing the stack operations (memory).*

*Define the following configurations:*

1. *$f_{\theta_c}^{frozen}$: The RNN/controller $\theta_c$ is frozen, and only the memory $\theta_m$ is trainable.*

2. *$f_\theta^{unstable}$: Both $\theta_c$ and $\theta_m$ are trainable, but the model is unstable.*

3. *$f_\theta^{frozen}$: Both the RNN/controller $\theta_c$ and memory $\theta_m$ are frozen (i.e., the system is untrained or random).*

*Let the loss function be $\text{Loss}(f_\theta(x), y)$ for input sequence $x \in \Sigma^*$ and target $y$. The following results hold:*

*1. **Partially Frozen Model Can Outperform an Unstable Model**:*

*There exists a range $T \in [T_{low}, T_{high}]$ such that:*

$$\mathbb{E}\left[ Loss(f_{\theta_c}^{frozen}(x), y) \right] < \mathbb{E}\left[ Loss(f_\theta^{unstable}(x), y) \right],$$

*for $x$ where $|x| \in [T_{low}, T_{high}]$.*

*2. **Error Growth Bound for Unstable Models**:*

*For the fully trained unstable model, the error can exhibit super-linear growth, which can be bounded as:*

$$Loss(f_\theta^{unstable}(x), y) \leq a \cdot |x|^b + c,$$

*where $b > 1$ indicates super-linear growth, and $a, c > 0$ are constants dependent on the degree of instability.*

*3. **Performance Bound for Partially Frozen Model**:*

*The error for a frozen RNN with trainable memory is bounded by:*

$$Loss(f_{\theta_c}^{frozen}(x), y) \leq a' \cdot |x| + c',$$

*where $a', c' > 0$ are constants, and $a'$ is typically smaller than the corresponding constant $a$ for the unstable model. This linear growth ensures better performance for a stable stack-augmented model even if the controller is frozen.*

We provide detailed proof in appendix and prove each conditions highlighted above.

We then examine error bounds in unstable models, demonstrating how these models might initially perform well on slightly longer sequences but eventually deteriorate as sequence length increases. This degradation is particularly relevant when assessing models designed to recognize languages with nested structures. The theorems characterize the error growth in unstable systems, providing insight into the rapid error accumulation and eventual convergence to random guessing.

**Theorem 0.7** (Error Bounds of an Unstable System). *Let $f_\theta$ be a stack-augmented Recurrent Neural Network (RNN) parameterized by $\theta = (\theta_c, \theta_m)$, where:*

- *$\theta_c$ represents the parameters of the RNN/controller,*

- *$\theta_m$ represents the parameters governing the stack operations (memory).*

*The system is considered unstable if it fails to maintain consistent performance as sequence length increases. Let the loss function be $Loss(f_\theta(x), y)$ for input sequence $x \in \Sigma^*$ and target $y$. Then the following results hold:*

*1. **Initial Performance on Slightly Longer Sequences**:*

*There exists a range of sequence lengths $|x| \in [T_{train}, T_{train} + \Delta T]$, where $\Delta T$ is small, such that the model may exhibit near-perfect performance (e.g., low error):*

$$Loss(f_\theta(x), y) \approx 0 \quad for \ |x| \in [T_{train}, T_{train} + \Delta T].$$

*2. **Error Growth Beyond the Initial Range**:*

*As the sequence length continues to increase beyond $T_{train} + \Delta T$, the error grows rapidly, leading to:*

$$Loss(f_\theta(x), y) \leq a \cdot |x|^b + c \quad for \ |x| > T_{train} + \Delta T,$$

*where $a > 0$, $b > 1$, and $c > 0$ indicate super-linear error growth.*

*3. **Convergence to Random Guessing or Worse**:*

*As sequence length $|x|$ continues to grow, the model's performance deteriorates further, ultimately converging to random guessing. The expected loss can be bounded as:*

$$\lim_{|x| \to \infty} Loss(f_\theta(x), y) \geq Loss_{random} = \frac{1}{k},$$

*where $k$ is the number of classes. In some cases, the error may exceed this bound due to the instability:*

$$\lim_{|x| \to \infty} Loss(f_\theta(x), y) \geq \frac{1}{k} + \epsilon \quad for \ some \ small \ \epsilon > 0.$$

The detailed proof is in the appendix. In the next section, we conduct a series of experiments showing that our methodology connects these theoretical results to practical considerations.

## Experimental Setup

We evaluate various RNN architectures on non-context-free languages (non-CFLs) and natural language modeling tasks. For the Penn Treebank (PTB) dataset, we test five models: a standard LSTM [1], two stack-augmented RNNs [7], and two nondeterministic stack-augmented RNNs [12] (see appendix for details). For non-CFL tasks, we also evaluate one LSTM [1], three stack-augmented RNNs [7], and one nondeterministic stack-augmented RNN [12].

**Non-Context-Free Languages** We investigate how stack-augmented RNNs handle non-CFL phenomena by evaluating them on several complex language modeling tasks. Each of the non-CFLs in our study can be recognized by a real-time three-stack automaton. We examine the following seven tasks:

1. **Count-3**: The language $a^n b^n c^n \mid n \geq 0$, which involves counting, reversing, and copying strings with markers.

2. **Marked-Reverse-and-Copy**: The language $\{w \# w^R \# w \mid w \in \{0, 1\}^*\}$, which requires reversing and copying a string with explicit markers.

| model_name | bin0 | bin1 | bin2 |
|---|---|---|---|
| lstm (n) | **0.98** | 0.87 | 0.75 |
| lstm (c) | 0.96 | **0.98** | **0.99** |
| jm-hidden (n) | 0.99 | 0.94 | 0.81 |
| jm-hidden (m) | 0.99 | **0.99** | **0.97** |
| rns-3-3 (n) | 0.98 | **0.94** | **0.85** |
| rns-3-3 (m) | 0.98 | 0.85 | 0.72 |

(a) count-3

| model_name | bin0 | bin1 | bin2 |
|---|---|---|---|
| lstm (n) | **0.76** | 0.49 | 0.48 |
| lstm (c) | 0.49 | **0.50** | **0.50** |
| jm-hidden (n) | **0.79** | 0.48 | 0.48 |
| jm-hidden (m) | 0.78 | **0.56** | **0.50** |
| rns-3-3 (n) | **0.82** | **0.75** | **0.67** |
| rns-3-3 (m) | 0.78 | 0.47 | 0.41 |

(b) marked-reverse-and-copy

| model_name | bin0 | bin1 | bin2 |
|---|---|---|---|
| lstm (n) | **0.73** | 0.61 | 0.61 |
| lstm (c) | 0.64 | **0.65** | **0.66** |
| jm-hidden (n) | **0.80** | 0.61 | 0.57 |
| jm-hidden (m) | 0.76 | 0.60 | 0.54 |
| rns-3-3 (n) | **0.81** | **0.69** | **0.55** |
| rns-3-3 (m) | 0.74 | 0.59 | 0.53 |

(c) count-and-copy

| model_name | bin0 | bin1 | bin2 |
|---|---|---|---|
| lstm (n) | **0.65** | 0.49 | 0.43 |
| lstm (c) | 0.49 | **0.49** | **0.50** |
| jm-hidden (n) | **0.68** | 0.42 | 0.43 |
| jm-hidden (m) | 0.63 | **0.47** | **0.50** |
| rns-3-3 (n) | **0.73** | 0.50 | 0.32 |
| rns-3-3 (m) | 0.57 | **0.38** | **0.32** |

(d) marked-copy

| model_name | bin0 | bin1 | bin2 |
|---|---|---|---|
| lstm (n) | **0.68** | **0.56** | **0.52** |
| lstm (c) | 0.49 | 0.49 | 0.50 |
| jm-10 (n) | **0.70** | **0.53** | **0.43** |
| jm-10 (m) | 0.68 | 0.51 | 0.39 |
| rns-3-3 (n) | **0.69** | 0.52 | 0.41 |
| rns-3-3 (m) | 0.65 | **0.52** | **0.49** |

(e) unmarked-copy-different-alphabets

| model_name | bin0 | bin1 | bin2 |
|---|---|---|---|
| lstm (n) | **0.67** | **0.54** | **0.51** |
| lstm (c) | 0.51 | 0.50 | 0.50 |
| jm-10 (n) | **0.70** | **0.56** | **0.51** |
| jm-10 (m) | 0.66 | 0.54 | 0.50 |
| rns-3-3 (n) | **0.69** | 0.53 | 0.50 |
| rns-3-3 (m) | 0.69 | **0.56** | **0.51** |

(f) unmarked-reverse-and-copy

Table 1: Test accuracy on the six context-sensitive languages of the best of 10 random restarts for each architecture. Bin0, bin1, and bin2 contain test set in the range $[40 - 100]$, $[100 - 200]$ and $[200 - 400]$, respectively, where $n$ represents models where all parameters are trained; whereas $m$ represents models where only memory is trained and rest all parameters are kept random

| model_name | bin0 | bin1 | bin2 |
|---|---|---|---|
| lstm (n) | **0.60** | **0.54** | **0.51** |
| lstm (c) | 0.49 | 0.50 | 0.50 |
| jm-hidden (n) | **0.63** | **0.55** | **0.51** |
| jm-hidden (m) | 0.62 | 0.54 | 0.50 |
| rns-3-3 (n) | 0.56 | 0.50 | 0.50 |
| rns-3-3 (m) | **0.60** | **0.53** | **0.50** |

Table 2: Test accuracy on the unmarked-copy CFL of the best of 10 random restarts for each architecture over bin0, bin1, and bin2 respectively

3. **Count-and-Copy**: The language $\{w\#^n w \mid w \in 0, 1\}$, which requires counting and copying strings separated by marked divisions.

4. **Marked-Copy**: The language $\{w\#w \mid w \in 0, 1^*\}$, which involves a simple marked copying operation.

5. **Unmarked-Copy-Different-Alphabets**: The language $\{ww' \mid w \in 0, 1, w' = \phi(w)\}$, where $\phi$ is a homomorphism defined as $\phi(0) = 2, \phi(1) = 3$.

6. **Unmarked-Reverse-and-Copy**: The language $\{ww^R w \mid w \in 0, 1\}$, which involves reversing and copying without explicit markers.

7. **Unmarked-Copy**: The language $\{ww \mid w \in 0, 1^*\}$, which requires unmarked copying of sequences.

For consistency, we followed the experimental framework and hyperparameters set by prior work [12]. The training and validation sets were sampled from sequences in the length range of $[40, 80]$, while the test sets were split into three bins based on sequence length: bin 0 ($[40, 100]$), bin 1 ($[100, 200]$), and bin 2 ($[200, 400]$).

**Natural Language Modeling** In addition to non-CFL tasks, we evaluated the stack-augmented RNNs on nat-

ural language modeling using the Penn Treebank (PTB) dataset, preprocessed as in Mikolov et al. (2011). The hyperparameters used for these experiments were consistent with those from prior work [12].

| Model Name | Test PPL (n) | Test PPL (m) |
|---|---|---|
| lstm-256 | **119.8** | 129.8 |
| jm-hidden-247 | 126.8 | 125.1 |
| jm-learned-22 | 125.9 | 124.2 |
| rns-4-5 | 126.5 | **120.5** |
| vrns-3-3-5 | 123.5 | 126.0 |

Table 3: Test and validation perplexity on the Penn Treebank of the best of 10 random restarts for each architecture, where $n$ represents models with all trainable parameters, whereas $m$ represents models where only memory is trained.

## Discussion and Conclusion

Our experiments on seven CGL benchmarks highlight the critical role of stability in understanding the learnability of stack-augmented RNNs. As shown in Table 1 and 2, both fully trained models and those with only memory trained perform comparably on short sequences. However, when tested on longer sequences (bin2), most models converge to random guessing, underscoring the instability introduced by increased sequence complexity. For instance, in the *count-3* task, fully trained models experience a significant accuracy drop—from 98% to 99% down to 75% to 85%. In contrast, freezing the controller during training allows models like LSTM and jm-hidden to maintain around 96% accuracy, while fixing both the controller and memory enables four out of five models to preserve their initial performance. This demonstrates that excessive parameter tuning can destabilize learning rather than enhance it. A similar trend is observed in the *count-and-copy* task, where models with fixed components maintain stable accuracy near their initial 64%, while fully trained models suffer a drop from 73% to 81% down to 55% to 61%. The PTB dataset also shows that models with only memory trained outperform fully trained counterparts, reinforcing that over-parameterization during training introduces noise and instability (Table 3). These results highlight that stability is essential for effective learning. Models exhibiting stable performance across varying sequence lengths are more likely to generalize well, while instability often leads to poor long-term retention and performance degradation. Additional analysis are shown in appendix.

## Conclusion.

Our theoretical analysis demonstrates that unstable systems tend to converge to random guessing, or even worse, under certain conditions. We validated these findings through a series of experiments that consistently showed this behavior in practice. These results highlight that understanding stability and rigorously testing models on longer sequences are crucial for achieving better learnability. In scenarios where longer sequences are not available, the error bounds of the system ideally should be better than random guessing or models with partially frozen parameters. Our study emphasizes that selectively freezing components during training can lead to more stable behavior and improved generalization, especially for longer sequences but is still farway compared to stable models. Thus evaluating stability across varying sequence lengths provides key insights into a model's ability to maintain robust performance, making stability assessments a critical factor in building more reliable and learnable sequence models.

## References

[1] F. A. Gers and E. Schmidhuber, "Lstm recurrent networks learn simple context-free and context-sensitive languages," *IEEE Transactions on Neural Networks*, vol. 12, pp. 1333–1340, Nov 2001.

[2] W. Merrill, G. Weiss, Y. Goldberg, R. Schwartz, N. A. Smith, and E. Yahav, "A formal hierarchy of rnn architectures," *arXiv preprint arXiv:2004.08500*, 2020.

[3] A. Mali, A. Ororbia, D. Kifer, and L. Giles, "On the computational complexity and formal hierarchy of second order recurrent neural networks," *arXiv preprint arXiv:2309.14691*, 2023.

[4] A. A. Mali, A. G. Ororbia II, and C. L. Giles, "A neural state pushdown automata," *IEEE Transactions on Artificial Intelligence*, vol. 1, no. 3, pp. 193–205, 2020.

[5] J. Stogin, A. Mali, and C. L. Giles, "A provably stable neural network turing machine with finite precision and time," *Information Sciences*, vol. 658, p. 120034, 2024.

[6] A. Graves, G. Wayne, M. Reynolds, T. Harley, I. Danihelka, A. Grabska-Barwińska, S. G. Colmenarejo, E. Grefenstette, T. Ramalho, J. Agapiou, *et al.*, "Hybrid computing using a neural network with dynamic external memory," *Nature*, vol. 538, no. 7626, p. 471, 2016.

[7] A. Joulin and T. Mikolov, "Inferring algorithmic patterns with stack-augmented recurrent nets," in *Advances in neural information processing systems*, pp. 190–198, 2015.

[8] E. Grefenstette, K. M. Hermann, M. Suleyman, and P. Blunsom, "Learning to transduce with unbounded memory," in *Advances in neural information processing systems*, pp. 1828–1836, 2015.

[9] B. DuSell and D. Chiang, "Learning context-free languages with nondeterministic stack rnns," *arXiv preprint arXiv:2010.04674*, 2020.

[10] A. Mali, A. G. Ororbia, D. Kifer, and C. L. Giles, "Recognizing and verifying mathematical equations using multiplicative differential neural units," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 35, pp. 5006–5015, 2021.

[11] J. E. Hopcroft, R. Motwani, and J. D. Ullman, "Introduction to automata theory, languages, and computation," *Acm Sigact News*, vol. 32, no. 1, pp. 60–65, 2001.

[12] B. DuSell and D. Chiang, "The surprising computational power of nondeterministic stack RNNs," in *The Eleventh International Conference on Learning Representations*, 2023.

[13] S. Hölldobler, Y. Kalinke, and H. Lehmann, "Designing a counter: Another case study of dynamics and activation landscapes in recurrent networks," in *KI-97: Advances in Artificial Intelligence: 21st Annual German Conference on Artificial Intelligence Freiburg, Germany, September 9–12, 1997 Proceedings 21*, pp. 313–324, Springer, 1997.

[14] M. Steijvers and P. Grünwald, "A recurrent network that performs a context-sensitive prediction task," in *Proceedings of the Eighteenth Annual Conference of the Cognitive Science Society*, pp. 335–339, Routledge, 2019.

[15] N. Skachkova, T. A. Trost, and D. Klakow, "Closing brackets with recurrent neural networks," in *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pp. 232–239, 2018.

[16] L. Sennhauser and R. C. Berwick, "Evaluating the ability of lstms to learn context-free grammars," *arXiv preprint arXiv:1811.02611*, 2018.

[17] M. Bodén and J. Wiles, "Context-free and context-sensitive dynamics in recurrent neural networks," *Connection Science*, vol. 12, no. 3-4, pp. 197–210, 2000.

[18] F. A. Gers and E. Schmidhuber, "Lstm recurrent networks learn simple context-free and context-sensitive languages," *IEEE transactions on neural networks*, vol. 12, no. 6, pp. 1333–1340, 2001.

[19] C. Anil, Y. Wu, A. Andreassen, A. Lewkowycz, V. Misra, V. Ramasesh, A. Slone, G. Gur-Ari, E. Dyer, and B. Neyshabur, "Exploring length generalization in large language models," *Advances in Neural Information Processing Systems*, vol. 35, pp. 38546–38556, 2022.

[20] N. Dave, D. Kifer, C. L. Giles, and A. Mali, "Investigating symbolic capabilities of large language models," *IJCAI 2024 Workshop on Logical Foundations of Neuro-Symbolic AI*, 2024.

[21] N. Lan, M. Geyer, E. Chemla, and R. Katzir, "Minimum description length recurrent neural networks," *Transactions of the Association for Computational Linguistics*, vol. 10, pp. 785–799, 2022.

[22] A. Mali, A. Ororbia, D. Kifer, and L. Giles, "Investigating backpropagation alternatives when learning to dynamically count with recurrent neural networks," in *International Conference on Grammatical Inference*, pp. 154–175, PMLR, 2021.

[23] Y. Chen, S. Gilroy, A. Maletti, J. May, and K. Knight, "Recurrent neural networks as weighted language recognizers," *arXiv preprint arXiv:1711.05408*, 2017.

[24] J. Pérez, J. Marinković, and P. Barceló, "On the turing completeness of modern neural network architectures," *arXiv preprint arXiv:1901.03429*, 2019.

[25] J. Ackerman and G. Cybenko, "A survey of neural networks and formal languages," *arXiv preprint arXiv:2006.01338*, 2020.

[26] S. Bhattamishra, K. Ahuja, and N. Goyal, "On the ability and limitations of transformers to recognize formal languages," 2020.

[27] A. Mali, A. Ororbia, D. Kifer, and L. Giles, "On the tensor representation and algebraic homomorphism of the neural state turing machine," *https://arxiv.org/pdf/2309.14690v1.pdf*.

[28] I. Danihelka, G. Wayne, B. Uria, N. Kalchbrenner, and A. Graves, "Associative long short-term memory," in *International conference on machine learning*, pp. 1986–1994, PMLR, 2016.

[29] K. Kurach, M. Andrychowicz, and I. Sutskever, "Neural random-access machines," *arXiv preprint arXiv:1511.06392*, 2015.

[30] A. Graves, G. Wayne, and I. Danihelka, "Neural turing machines," *arXiv preprint arXiv:1410.5401*, 2014.

## Appendix A: Related Work

The capability of Recurrent Neural Networks (RNNs), Gated Recurrent Units (GRUs), and Long Short-Term Memory networks (LSTMs) to learn formal languages has been well-studied. While these architectures perform well on simple languages, such as basic counting tasks and Dyck languages [13, 14, 15], they struggle with more complex languages, especially when generalizing to longer sequences [16]. The lack of external memory structures limits these models' ability to handle advanced languages, confining their generalization to sequence lengths close to those seen during training [17, 18]. While some recent work has explored length generalization in synthetic reasoning tasks with large pretrained models [19, 20], these studies are often disconnected from the more rigorous demands of formal language learning. On other hand several studies have shown that the choice of objective functions [21] and learning algorithms [22] significantly affects RNNs' ability to stably learn complex grammars. For instance, [21] demonstrated that specialized loss functions, such as minimum description length, lead to more stable convergence and better generalization on formal language tasks.

Although RNNs are theoretically Turing complete, practical constraints such as finite precision and limited recurrent steps diminish their real-world expressivity [23, 24]. Under such constraints, standard RNNs and GRUs are restricted to regular languages, while LSTMs can approximate context-free languages by simulating k-counter mechanisms [25, 26, 2]. Tensor RNNs [3, 27] have extended these capabilities, but their higher computational costs make them challenging to deploy.

Memory-augmented architectures have been proposed to overcome these limitations by integrating external structures such as stacks [5, 7, 8], random access memory [28, 29], and memory matrices [30, 6]. These models have shown success in recognizing more complex languages and executing tasks like string copying and reversal. However, most research has focused on theoretical expressivity and evaluations on short sequences. Consequently, many studies fail to test these models on longer sequences, where stability issues become more apparent [20]. Empirical studies have begun to emphasize the need for testing on longer sequences [4, 10, 27], and recent theoretical work by Stogin [5] has identified stable differentiable memory structures. Yet, these insights are not comprehensive and lack a systematic analysis of stability.

The primary gap in existing research is the absence of focused investigations into the stability of memory-augmented models, particularly when processing sequences significantly longer than those seen during training. Stability, in this context, refers to a model's ability to maintain consistent performance across varying sequence lengths without degradation in accuracy or memory manipulation. Most studies evaluate models on sequences similar in length to the training data, thereby overlooking the challenges posed by longer inputs. As a result, while these models often achieve high accuracy within the training range, they frequently exhibit instability when exposed to longer sequences.

This study addresses these gaps by examining the stability of stack-augmented RNNs across varying sequence lengths. We explore how different configurations, such as freezing the RNN controller while keeping the memory trainable, impact stability. We derive theoretical conditions for stable performance and identify when models become unstable. Our work is among the first to systematically investigate the stability of memory-augmented neural networks from both theoretical and empirical perspectives, providing insights into factors that influence learnability and generalization across longer and more complex sequences.

## Appendix B: Additional Results

Figure 1-3 and Table 4-7 provide comprehensive overview of our hypothesis. In our experiments with the non-context-free language task *count-3*, a notable divergence in performance emerged between fully trained models and those with selective components frozen. During evaluation on sequences within the training range (40 to 100 tokens), all models performed comparably. However, as sequence lengths extended beyond this range (100-200 and 200-400 tokens), fully trained models exhibited significant performance degradation, with accuracy decreasing to between 75% and 85%. This decline illustrates the difficulty these models face in generalizing to longer sequences, where learned dependencies may not scale effectively.

Conversely, models with a frozen controller demonstrated remarkable stability, maintaining near-perfect accuracy (up to 99%) even on the longest test sequences. This suggests that freezing the controller helps mitigate instability introduced by continuous parameter updates, thereby improving the model's ability to generalize over extended input lengths.

Further analysis of the model *jm-hidden*, which leverages a superposition stack RNN and offloads the controller's hidden state to external memory, reinforces these observations. When fully trained (*jm-hidden none*), the model's predictions begin to diverge from the expected outputs as sequence lengths increase, particularly beyond 200 tokens, highlighting instability in handling extended input sequences. However, when the memory component is frozen (*jm-hidden m*), the model exhibits stable performance, closely aligning with the expected outputs across all tested sequence lengths, including the longest sequences.

This stability trend was consistently observed in our experiments on the Penn Treebank (PTB) dataset for language modeling. Models such as *jm-hidden-247*, *jm-learned-22*, and *rns-4-5* with frozen memory components not only retained performance but, in some cases, outperformed their fully trained counterparts, achieving superior perplexity scores (126.8 $\rightarrow$ 125.1, 125.9 $\rightarrow$ 124.2, 126.5 $\rightarrow$ 120.5, respectively). These results underscore the critical role of stability in enhancing the generalization capabilities of neural architectures, particularly in tasks involving long sequences. Freezing specific components, such as memory or the controller, leads to more reliable and robust performance, highlighting the necessity of incorporating stability considerations into neural network design.
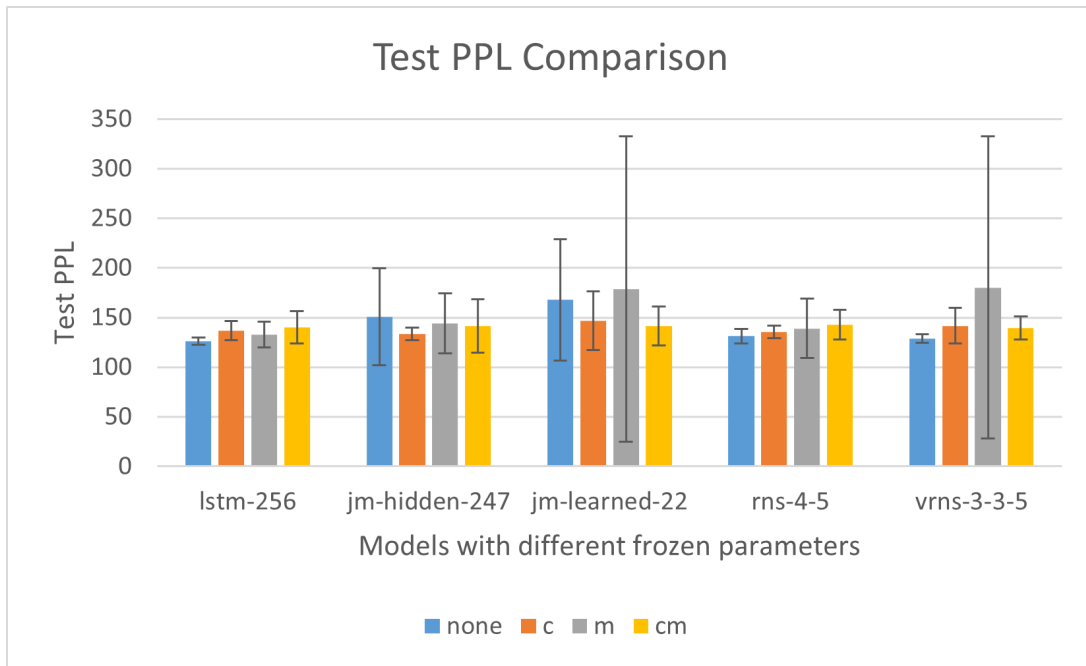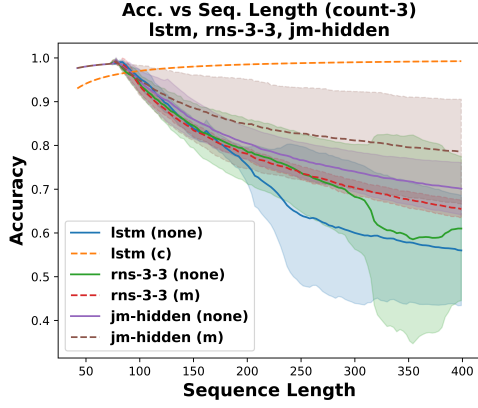


Figure 1: Performance of various models using 4 configuration (none = fully trained model, m = only memory is trained, c = only controller is trained and cm = controller and memory are frozen and only classifier is trainable. We report performance on language modeling task and report perplexity (PPL) on Penn tree bank (PTB) dataset.

## Appendix C: Evaluating Stability in Practice

In this section we provide various ways the stability in stack-augmented Recurrent Neural Networks (RNNs) can be practically assessed by focusing on the model's error behavior, generalization to longer sequences, and resilience to perturbations. The following conditions are essential for determining whether the system remains stable across varying input scenarios.

### Error Bounds Across Sequence Lengths

A stable system should maintain a bounded error for sequences of varying lengths. Let the loss function be $\text{Loss}(f_\theta(x), y)$ for input $x$ and target $y$. Ideally, the error should satisfy:

(a) count-3

(b) marked-reverse-and-copy

(c) count-and-copy

(d) marked-copy

(e) unmarked-copy-different-alphabets

(f) unmarked-reverse-and-copy

(g) unmarked-copy

15

Figure 2: Performance of top 3 models on test sets across 7 context free languages, when models are fully trained (none) and when only memory (m) is frozen.

(a) count-3

(b) marked-reverse-and-copy

(c) count-and-copy

(d) marked-copy

(e) unmarked-copy-different-alphabets

(f) unmarked-reverse-and-copy

(g) unmarked-copy

16

Figure 3: Performance of top 2 memory-augmented models on test sets across 7 context free languages, when models are fully trained (none) and when only memory (m) is frozen.

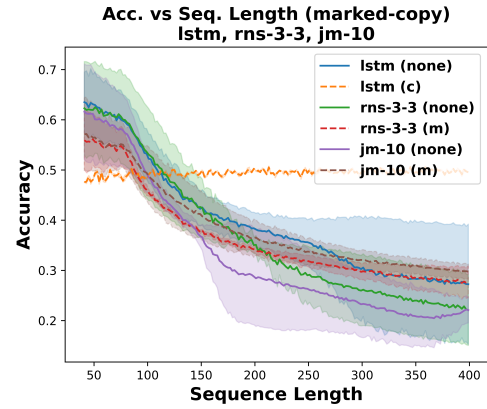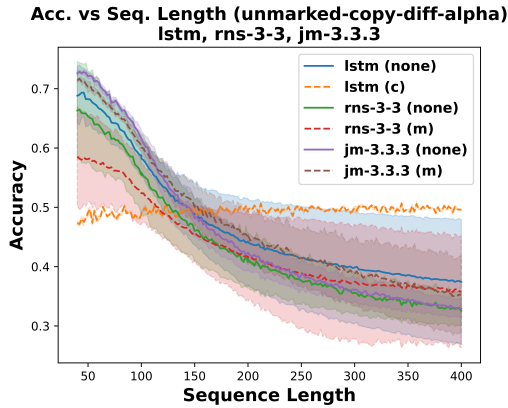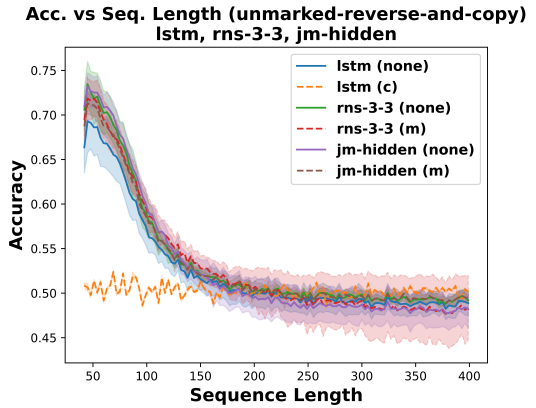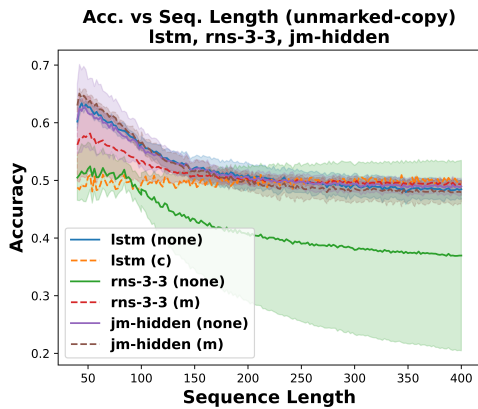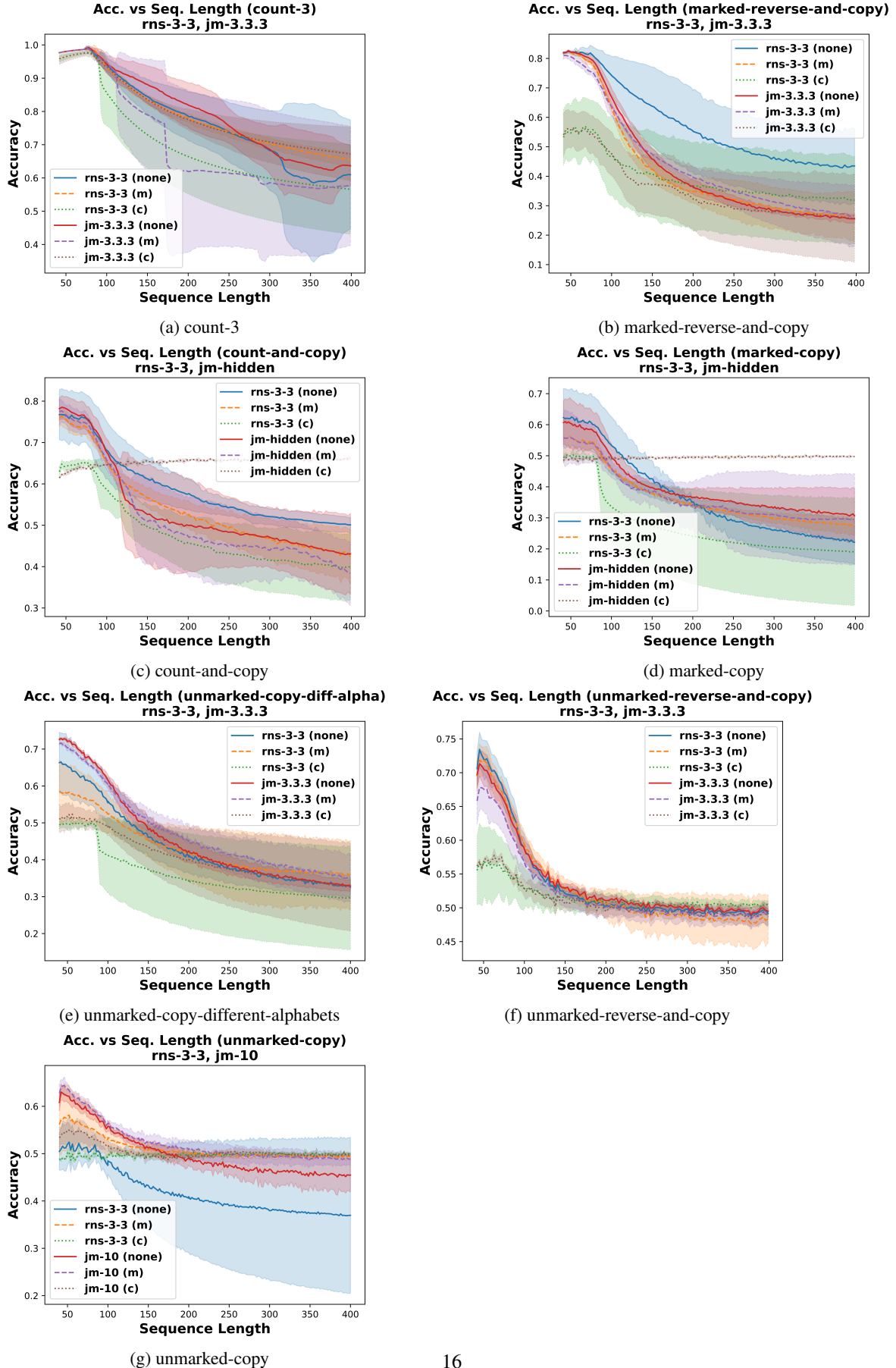| Task | Model | Metrics | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | | Val PPL | Bin0 PPL | Bin1 PPL | Bin2 PPL | Bin0 Acc | Bin1 Acc | Bin2 Acc |
| count-3 | lstm | **1.04** | 1.06 | 1.54 | 1772.19 | 0.98 | 0.87 | 0.75 |
| | jm-10 | **1.04** | **1.05** | 1.18 | 1.67 | **0.99** | 0.90 | 0.76 |
| | jm-3.3.3 | **1.04** | **1.05** | 1.42 | 11.23 | **0.99** | 0.93 | 0.78 |
| | jm-hidden | **1.04** | **1.05** | **1.14** | **1.22** | **0.99** | **0.94** | 0.81 |
| | rns-3-3 | **1.04** | 1.06 | 1.35 | 155.42 | 0.98 | **0.94** | **0.85** |
| marked-reverse-and-copy | lstm | 1.40 | 1.56 | 10.65 | **10.57** | 0.76 | 0.49 | 0.48 |
| | jm-10 | 1.33 | 1.42 | 3.55 | 182.41 | 0.80 | 0.50 | 0.43 |
| | jm-3.3.3 | 1.33 | 1.40 | 11.46 | 1060.93 | 0.80 | 0.50 | 0.31 |
| | jm-hidden | 1.32 | 1.49 | 103.81 | 12.71 | 0.79 | 0.48 | 0.48 |
| | rns-3-3 | **1.30** | **1.33** | **1.68** | 18.47 | **0.82** | **0.75** | **0.67** |
| count-and-copy | lstm | 1.48 | 1.57 | 13.71 | 18.75 | 0.73 | 0.61 | 0.61 |
| | jm-10 | 1.39 | 1.61 | 368.71 | 10.42 | 0.76 | 0.60 | 0.59 |
| | jm-3.3.3 | 1.37 | 1.56 | **1.97** | 4.93 | 0.76 | 0.64 | **0.62** |
| | jm-hidden | 1.34 | 1.41 | 13.04 | **3.53** | 0.80 | 0.61 | 0.57 |
| | rns-3-3 | **1.31** | **1.38** | 2.87 | 19.00 | **0.81** | **0.69** | 0.55 |
| marked-copy | lstm | 1.66 | 1.81 | 3.11 | 17.58 | 0.65 | 0.49 | 0.43 |
| | jm-10 | 1.63 | 1.77 | 14.13 | 16.13 | 0.66 | 0.40 | 0.30 |
| | jm-3.3.3 | 1.68 | 1.77 | **2.55** | **2.79** | 0.65 | **0.53** | **0.50** |
| | jm-hidden | 1.59 | 1.74 | 2.62 | 15.53 | 0.68 | 0.42 | 0.43 |
| | rns-3-3 | **1.47** | **1.51** | 11.05 | 2.82 | **0.73** | 0.50 | 0.32 |
| unmarked-copy-diff-alpha | lstm | 1.61 | 1.68 | 2.66 | 17.51 | 0.68 | **0.56** | **0.52** |
| | jm-10 | **1.56** | **1.65** | 3.10 | 12.99 | **0.70** | 0.53 | 0.43 |
| | jm-3.3.3 | 1.57 | **1.65** | 3.00 | 14.73 | 0.68 | 0.52 | 0.39 |
| | jm-hidden | 1.58 | 1.70 | 2.40 | **10.02** | 0.68 | 0.48 | 0.46 |
| | rns-3-3 | **1.56** | 1.66 | **2.20** | 13.60 | 0.69 | 0.52 | 0.41 |
| unmarked-reverse-and-copy | lstm | 1.66 | 1.84 | 3.32 | 3.93 | 0.67 | 0.54 | **0.51** |
| | jm-10 | **1.60** | **1.78** | 3.65 | 10.15 | **0.70** | **0.56** | **0.51** |
| | jm-3.3.3 | 1.64 | 1.83 | **3.24** | **3.66** | 0.68 | 0.54 | 0.50 |
| | jm-hidden | **1.60** | 1.81 | 4.21 | 11.75 | 0.69 | 0.54 | 0.50 |
| | rns-3-3 | **1.60** | **1.78** | 3.65 | 13.46 | 0.69 | 0.53 | 0.50 |
| unmarked-copy | lstm | 1.85 | 1.92 | 2.61 | 3.19 | 0.60 | 0.54 | 0.51 |
| | jm-10 | 1.82 | 1.91 | 2.49 | 2.69 | 0.61 | 0.53 | **0.50** |
| | jm-3.3.3 | 1.83 | 1.92 | 2.20 | 2.25 | 0.61 | 0.53 | **0.50** |
| | jm-hidden | **1.76** | **1.86** | **2.16** | **2.18** | 0.63 | 0.55 | 0.51 |
| | rns-3-3 | 1.94 | 2.06 | 2.19 | 2.23 | **0.56** | **0.50** | **0.50** |

Table 4: Performance of different models on context free languages when all models components are fully trained (none). We test all models on 3 independent test set, slightly long test set bin 0 (N = $[40 - 100]$), mid range test (N=$[100 - 200]$)set bin1 and bin2 (N = $[200 - 400]$

$$|\text{Loss}(f_\theta(x), y)| \leq C \quad \forall x \in \Sigma^* \text{ where } T_{\min} \leq |x| \leq T_{\max},$$

where $T_{\min}$ and $T_{\max}$ define the minimum and maximum sequence lengths tested. The goal is to ensure that the error remains within a predefined bound $C > 0$ across sequences of different lengths.

**Consistency of Error Across Sequence Lengths**

The error should be consistent across sequences of varying lengths. Variance in error as a function of sequence length should be minimal. Formally, the stability criterion is:

$$\text{Var}\left(\text{Loss}(f_\theta(x), y)\right) \text{ should be minimized across varying } |x|.$$

Low variance in error indicates that the model's behavior is predictable and stable regardless of the input sequence length.

17

| Task | Model | Metrics | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | | Val PPL | Bin0 PPL | Bin1 PPL | Bin2 PPL | Bin0 Acc | Bin1 Acc | Bin2 Acc |
| **count-3** | lstm | 1.19 | 1.17 | **1.12** | **1.10** | 0.96 | **0.98** | 0.99 |
| | jm-10 | 1.10 | 1.12 | 1.32 | 1.43 | 0.97 | 0.90 | 0.79 |
| | jm-3.3.3 | **1.05** | **1.08** | 1.57 | 14.65 | **0.98** | 0.85 | 0.76 |
| | jm-hidden | 1.15 | 1.16 | 1.16 | 1.18 | 0.96 | **0.98** | 0.99 |
| | rns-3-3 | **1.05** | 1.21 | 136.65 | 126512.54 | 0.97 | 0.83 | **0.75** |
| **marked-reverse-and-copy** | lstm | 2.38 | 2.35 | 2.23 | **2.18** | 0.49 | 0.50 | 0.50 |
| | jm-10 | 2.17 | 2.27 | 2.45 | 2.39 | 0.50 | 0.48 | 0.49 |
| | jm-3.3.3 | 1.83 | 1.92 | 10.43 | 12.49 | 0.66 | 0.50 | 0.50 |
| | jm-hidden | 2.27 | 2.26 | 2.21 | 2.19 | 0.56 | 0.53 | 0.51 |
| | rns-3-3 | **1.55** | **1.58** | **1.95** | 2.29 | **0.75** | **0.65** | **0.55** |
| **count-and-copy** | lstm | 1.93 | 1.89 | **1.79** | **1.74** | 0.64 | **0.65** | 0.66 |
| | jm-10 | 1.70 | 1.79 | 15.26 | 10.83 | 0.63 | 0.50 | 0.42 |
| | jm-3.3.3 | 1.66 | 1.81 | 10.60 | 108.02 | 0.63 | 0.43 | 0.30 |
| | jm-hidden | 1.85 | 1.84 | **1.79** | **1.74** | 0.64 | **0.65** | 0.66 |
| | rns-3-3 | **1.65** | **1.71** | 1.95 | 13.67 | **0.65** | 0.57 | 0.53 |
| **marked-copy** | lstm | 2.30 | 2.27 | 2.16 | **2.12** | **0.49** | 0.49 | **0.50** |
| | jm-10 | 2.11 | 2.20 | 2.37 | 17.03 | **0.49** | 0.49 | **0.50** |
| | jm-3.3.3 | 2.13 | 2.18 | 2.37 | 2.37 | **0.49** | **0.50** | **0.50** |
| | jm-hidden | 2.24 | 2.22 | **2.15** | **2.12** | **0.49** | 0.49 | **0.50** |
| | rns-3-3 | **2.07** | **2.16** | 2.90 | 13.93 | **0.49** | 0.49 | **0.50** |
| **unmarked-copy-diff-alpha** | lstm | 2.28 | 2.25 | **2.16** | **2.12** | 0.49 | 0.49 | **0.50** |
| | jm-10 | 2.04 | 2.14 | 13.00 | 10846.00 | 0.55 | 0.46 | 0.36 |
| | jm-3.3.3 | **2.03** | **2.12** | 19.82 | 2.62 | **0.56** | 0.49 | **0.50** |
| | jm-hidden | 2.21 | 2.21 | 2.19 | 2.17 | 0.51 | **0.50** | **0.50** |
| | rns-3-3 | 2.08 | 11.02 | 2.31 | 181.99 | 0.50 | 0.49 | **0.50** |
| **unmarked-reverse-and-copy** | lstm | 2.15 | 2.13 | **2.08** | **2.06** | 0.51 | 0.50 | 0.50 |
| | jm-10 | 2.05 | 2.12 | 2.58 | 2.81 | 0.55 | 0.50 | 0.49 |
| | jm-3.3.3 | 2.04 | 2.10 | 2.47 | 2.71 | 0.56 | 0.52 | 0.51 |
| | jm-hidden | 2.10 | 2.11 | 2.12 | 2.13 | 0.51 | 0.51 | 0.50 |
| | rns-3-3 | **1.87** | **1.97** | 2.17 | 2.20 | **0.62** | **0.55** | **0.52** |
| **unmarked-copy** | lstm | 2.15 | 2.13 | **2.08** | **2.06** | 0.49 | **0.50** | **0.50** |
| | jm-10 | **2.06** | 2.11 | 2.22 | 2.29 | **0.55** | **0.50** | **0.50** |
| | jm-3.3.3 | 2.03 | **2.09** | 2.20 | 2.27 | **0.55** | **0.50** | **0.50** |
| | jm-hidden | 2.11 | 2.11 | 2.14 | 2.16 | 0.51 | **0.50** | **0.50** |
| | rns-3-3 | 2.11 | 2.12 | 2.16 | 2.19 | 0.49 | **0.50** | **0.50** |

Table 5: Performance of different models on context free languages when only parameter belonging to controller are trainable (mode = c). We test all models on 3 independent test set, slightly long test set bin 0 (N = $[40 - 100]$), mid range test (N=$[100 - 200]$)set bin1 and bin2 (N = $[200 - 400]$.

**Performance Degradation Over Long Sequences**

Stability is also determined by the model's ability to handle sequences significantly longer than those seen during training. Let $x_{\text{long}}$ and $x_{\text{short}}$ be sequences of longer and shorter lengths, respectively. Ideally, the model should satisfy:

$$\text{Loss}(f_\theta(x_{\text{long}}), y_{\text{long}}) \approx \text{Loss}(f_\theta(x_{\text{short}}), y_{\text{short}}).$$

Testing on sequences 2-4 times longer than those seen during training helps reveal whether the system can generalize to more complex scenarios without significant error growth.

**Generalization to Unseen Sequence Patterns**

The system should generalize to novel input sequences that conform to the language $L$ but are not explicitly encountered during training. The criterion is that the model's error on unseen patterns should be similar to that on known patterns:

$$|\text{Loss}(f_\theta(x_{\text{new}}), y_{\text{new}}) - \text{Loss}(f_\theta(x_{\text{known}}), y_{\text{known}})| \leq \delta,$$

| Task | Model | Metrics | | | | | | |
|------|-------|---------|---|---|---|---|---|---|
| | | Val PPL | Bin0 PPL | Bin1 PPL | Bin2 PPL | Bin0 Acc | Bin1 Acc | Bin2 Acc |
| **count-3** | lstm | **1.04** | **1.05** | 1.28 | 1.88 | **0.99** | 0.90 | 0.75 |
| | jm-10 | **1.04** | **1.05** | 1.21 | 1.70 | 0.98 | 0.86 | 0.76 |
| | jm-3.3.3 | **1.04** | 1.06 | 1.25 | **1.69** | 0.98 | 0.92 | 0.78 |
| | jm-hidden | **1.04** | **1.05** | **1.19** | 1.70 | **0.99** | **0.99** | **0.97** |
| | rns-3-3 | **1.04** | **1.05** | 1.36 | 2.87 | 0.98 | 0.85 | 0.72 |
| **marked-reverse-and-copy** | lstm | 1.38 | 1.55 | **11.64** | **3.56** | 0.77 | 0.48 | 0.48 |
| | jm-10 | 1.38 | 1.51 | 20.13 | 10.30 | 0.77 | 0.49 | 0.47 |
| | jm-3.3.3 | 1.41 | 1.54 | 12.47 | 11.30 | 0.76 | 0.49 | 0.46 |
| | jm-hidden | 1.37 | **1.49** | 134.47 | 13.62 | **0.78** | **0.56** | **0.50** |
| | rns-3-3 | **1.36** | 1.52 | 23.70 | 233.30 | **0.78** | 0.47 | 0.41 |
| **count-and-copy** | lstm | 1.43 | **1.50** | **1.93** | 16.86 | **0.76** | 0.63 | 0.54 |
| | jm-10 | 1.45 | 1.53 | 2.01 | **2.39** | **0.76** | **0.64** | **0.56** |
| | jm-3.3.3 | 1.49 | 1.56 | 2.41 | 1266.62 | 0.73 | 0.60 | 0.53 |
| | jm-hidden | **1.40** | 1.57 | 131.80 | 115.32 | **0.76** | 0.60 | 0.54 |
| | rns-3-3 | 1.43 | 1.58 | 16.84 | 175.87 | 0.74 | 0.59 | 0.53 |
| **marked-copy** | lstm | 1.65 | **1.77** | 2.73 | 10.88 | **0.67** | 0.49 | 0.42 |
| | jm-10 | 1.75 | 1.81 | 2.32 | 2.92 | 0.63 | **0.50** | 0.36 |
| | jm-3.3.3 | 1.76 | 1.85 | **2.28** | **2.26** | 0.62 | 0.49 | **0.50** |
| | jm-hidden | **1.63** | 1.89 | 159.06 | 103370.22 | 0.63 | 0.47 | **0.50** |
| | rns-3-3 | 1.79 | 2.02 | 18.16 | 2.82 | 0.57 | 0.38 | 0.32 |
| **unmarked-copy-diff-alpha** | lstm | 1.62 | **1.69** | **2.26** | 1103.65 | 0.67 | 0.54 | 0.47 |
| | jm-10 | **1.59** | 1.72 | 2.49 | 1242.00 | **0.68** | 0.51 | 0.39 |
| | jm-3.3.3 | 1.63 | 1.72 | 2.27 | 17.13 | 0.67 | **0.55** | 0.46 |
| | jm-hidden | 1.60 | 1.70 | 2.32 | 17.39 | 0.67 | 0.54 | **0.51** |
| | rns-3-3 | 1.71 | 1.82 | 13.09 | **13.72** | 0.65 | 0.52 | 0.49 |
| **unmarked-reverse-and-copy** | lstm | 1.64 | **1.81** | 3.12 | **10.05** | 0.68 | 0.54 | 0.50 |
| | jm-10 | 1.68 | 1.86 | **3.02** | 10.41 | 0.66 | 0.54 | 0.50 |
| | jm-3.3.3 | 1.68 | 1.87 | 3.21 | 11.97 | 0.67 | 0.53 | 0.50 |
| | jm-hidden | 1.65 | 1.83 | 3.45 | 14.59 | 0.67 | 0.53 | 0.50 |
| | rns-3-3 | **1.61** | 1.82 | 3.29 | 11.07 | **0.69** | **0.56** | **0.51** |
| **unmarked-copy** | lstm | 1.85 | 1.93 | **2.19** | **2.23** | 0.60 | **0.54** | **0.51** |
| | jm-10 | 1.82 | 1.92 | 2.44 | 2.67 | 0.61 | **0.54** | **0.51** |
| | jm-3.3.3 | 1.84 | 1.94 | 2.21 | 14.98 | 0.61 | 0.53 | 0.50 |
| | jm-hidden | **1.81** | **1.90** | 2.68 | 10.87 | **0.62** | **0.54** | 0.50 |
| | rns-3-3 | 1.85 | 1.94 | 2.27 | 2.38 | 0.60 | 0.53 | 0.50 |

Table 6: Performance of different models on context free languages when only parameter belonging to memory are trainable (mode = m). We test all models on 3 independent test set, slightly long test set bin 0 (N = $[40 - 100]$), mid range test (N=$[100 - 200]$)set bin1 and bin2 (N = $[200 - 400]$

where $x_{\text{new}}$ and $x_{\text{known}}$ represent novel and known sequences, respectively, and $\delta > 0$ is a small acceptable difference in error.

**Analyzing Error Growth**

The growth of error as sequence length increases is a critical measure of stability. For a stable system, error growth should be sub-linear or, at most, linear:

$$|\text{Loss}(f_\theta(x), y)| \leq C_1 \cdot |x| + C_2,$$

where $C_1$ and $C_2$ are constants. Super-linear growth (e.g., quadratic or exponential) indicates instability and should be avoided.

**Robustness to Perturbations**

A stable model should be resilient to small perturbations in the input. For a small perturbation $\epsilon$ added to the input sequence $x$, the error difference should be minimal:

$$|\text{Loss}(f_\theta(x + \epsilon), y) - \text{Loss}(f_\theta(x), y)| \leq \delta,$$

| Task | Model | Metrics | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | | Val PPL | Bin0 PPL | Bin1 PPL | Bin2 PPL | Bin0 Acc | Bin1 Acc | Bin2 Acc |
| **count-3** | lstm | 1.18 | 1.17 | **1.12** | **1.10** | **0.96** | **0.98** | **0.99** |
| | jm-10 | 1.17 | 1.17 | 1.20 | 1.33 | **0.96** | 0.96 | 0.83 |
| | jm-3.3.3 | 1.16 | 1.16 | 1.15 | 1.18 | **0.96** | **0.98** | **0.99** |
| | jm-hidden | **1.15** | **1.15** | **1.12** | 1.11 | **0.96** | **0.98** | **0.99** |
| | rns-3-3 | 1.19 | 1.17 | **1.12** | **1.10** | **0.96** | **0.98** | **0.99** |
| **marked-reverse-and-copy** | lstm | 2.38 | 2.35 | **2.23** | **2.18** | 0.49 | **0.50** | **0.50** |
| | jm-10 | **2.34** | **2.33** | 2.26 | 2.25 | **0.50** | **0.50** | **0.50** |
| | jm-3.3.3 | **2.34** | **2.33** | 2.28 | 2.30 | **0.50** | **0.50** | **0.50** |
| | jm-hidden | 2.37 | 2.34 | **2.23** | 2.19 | **0.50** | **0.50** | **0.50** |
| | rns-3-3 | 2.39 | 2.36 | **2.23** | **2.18** | 0.49 | **0.50** | **0.50** |
| **count-and-copy** | lstm | 1.93 | 1.89 | **1.79** | **1.74** | **0.64** | **0.65** | **0.66** |
| | jm-10 | **1.89** | 1.87 | 1.82 | 1.81 | **0.64** | **0.65** | **0.66** |
| | jm-3.3.3 | **1.89** | **1.86** | 1.81 | 1.80 | **0.64** | **0.65** | **0.66** |
| | jm-hidden | 1.90 | 1.87 | 1.80 | 1.76 | **0.64** | **0.65** | **0.66** |
| | rns-3-3 | 1.93 | 1.90 | **1.79** | **1.74** | **0.64** | **0.65** | **0.66** |
| **marked-copy** | lstm | 2.29 | 2.27 | **2.16** | **2.12** | **0.49** | **0.49** | **0.50** |
| | jm-10 | **2.26** | 2.25 | 2.19 | 2.18 | **0.49** | **0.49** | **0.50** |
| | jm-3.3.3 | **2.26** | **2.24** | 2.23 | 2.29 | **0.49** | **0.49** | **0.50** |
| | jm-hidden | 2.28 | 2.26 | 2.17 | **2.12** | **0.49** | **0.49** | **0.50** |
| | rns-3-3 | 2.29 | 2.26 | **2.16** | **2.12** | **0.49** | **0.49** | **0.50** |
| **unmarked-copy-diff-alpha** | lstm | 2.28 | 2.24 | **2.16** | **2.12** | 0.49 | 0.49 | **0.50** |
| | jm-10 | 2.24 | **2.22** | 2.19 | 2.20 | 0.49 | 0.49 | **0.50** |
| | jm-3.3.3 | **2.23** | **2.22** | 2.18 | 2.18 | 0.49 | **0.50** | **0.50** |
| | jm-hidden | 2.26 | 2.23 | **2.16** | **2.12** | **0.50** | **0.50** | **0.50** |
| | rns-3-3 | 2.28 | 2.24 | **2.16** | **2.12** | 0.49 | 0.49 | **0.50** |
| **unmarked-reverse-and-copy** | lstm | 2.15 | 2.13 | **2.08** | **2.06** | **0.51** | 0.50 | **0.50** |
| | jm-10 | **2.11** | **2.11** | 2.12 | 2.17 | **0.51** | 0.50 | **0.50** |
| | jm-3.3.3 | **2.11** | **2.11** | 2.10 | 2.11 | **0.51** | **0.51** | **0.50** |
| | jm-hidden | 2.13 | 2.12 | **2.08** | **2.06** | **0.51** | **0.51** | **0.50** |
| | rns-3-3 | 2.15 | 2.13 | **2.08** | **2.06** | **0.51** | 0.50 | **0.50** |
| **unmarked-copy** | lstm | 2.15 | 2.13 | **2.08** | **2.06** | 0.49 | **0.50** | **0.50** |
| | jm-10 | **2.12** | **2.11** | 2.13 | 2.22 | **0.50** | **0.50** | **0.50** |
| | jm-3.3.3 | **2.12** | **2.11** | 2.10 | 2.11 | **0.50** | **0.50** | **0.50** |
| | jm-hidden | 2.14 | 2.12 | **2.08** | **2.06** | **0.50** | **0.50** | **0.50** |
| | rns-3-3 | 2.15 | 2.13 | **2.08** | **2.06** | **0.50** | **0.50** | **0.50** |

Table 7: Performance of different models on context free languages when only classifier is trained; whereas parameter belonging to controller and memory are fixed (mode = cm). We test all models on 3 independent test set, slightly long test set bin 0 (N = $[40 - 100]$), mid range test (N=$[100 - 200]$)set bin1 and bin2 (N = $[200 - 400]$

where $\delta > 0$ is an acceptable margin for error change. High sensitivity to perturbations suggests instability.